# Sequencing Complex Mixtures of HIV-1 Genomes with Single-Base Resolution

Michael P. S. Brown[1], Melissa Laird[1], Lance Hepler[1], Yan Guo[1], Colleen Ludka[1], Ellen E. Paxinos[1], Richard Lempicki[2], Da Wei Huang[2], Cliff Lane[3], Hiromi Imamichi[3]

[1]Pacific Biosciences, Menlo Park, CA  [2]Leidos Biomedical Research, Frederick, MD  [3]NIAID, NIH, DHHS, Bethesda, MD
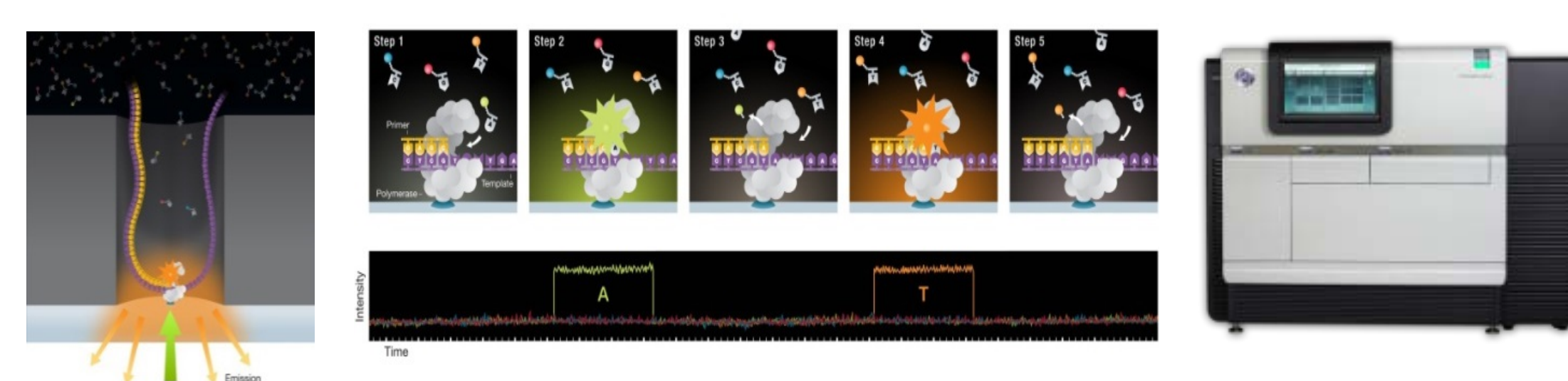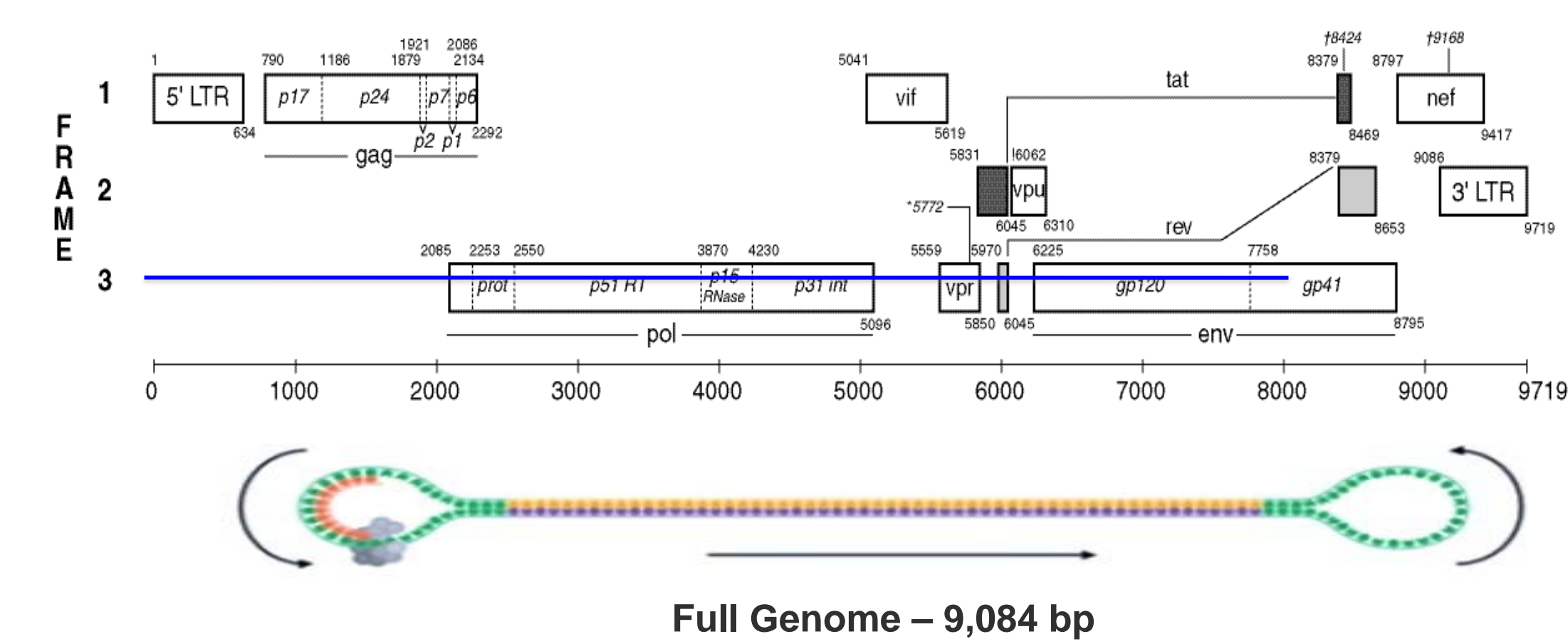
## Abstract

A large number of distinct HIV-1 genomes can be present in a single clinical sample from a patient chronically infected with HIV-1. We examined samples containing complex mixtures of near-full-length HIV-1 genomes. Single molecules were sequenced as near-full-length (9.6 kb) amplicons directly from PCR products without shearing.

Mathematical analysis techniques deconvolved the complex mixture of reads into estimates of distinct near-full-length viral genomes with their relative abundances. We correctly estimated the originating genomes to single-base resolution along with their relative abundances for mixtures where the truth was known exactly by independent sequencing methods. Correct estimates were made even when genomes diverged by a single base.

Minor abundances of 5% were reliably detected. SMRT® Sequencing data contained near-full-length continuous reads for each sample including some runs with greater than 10,000 near-full-length-genome reads in a three-hour collection time. SMRT Sequencing yields long-read sequencing results from individual DNA molecules with a rapid time-to-result. The single-molecule, full-length nature of the sequencing method allows us to estimate variant subspecies and relative abundances even from samples containing complex mixtures of genomes that differ by single bases.

These results open the possibility of cost-effective full-genome sequencing of HIV-1 in mixed populations for applications such as incorporated-HIV-1 screening. In screening, genomes can differ by one to many thousands of bases and the ability to measure them can help scientifically inform treatment strategies.

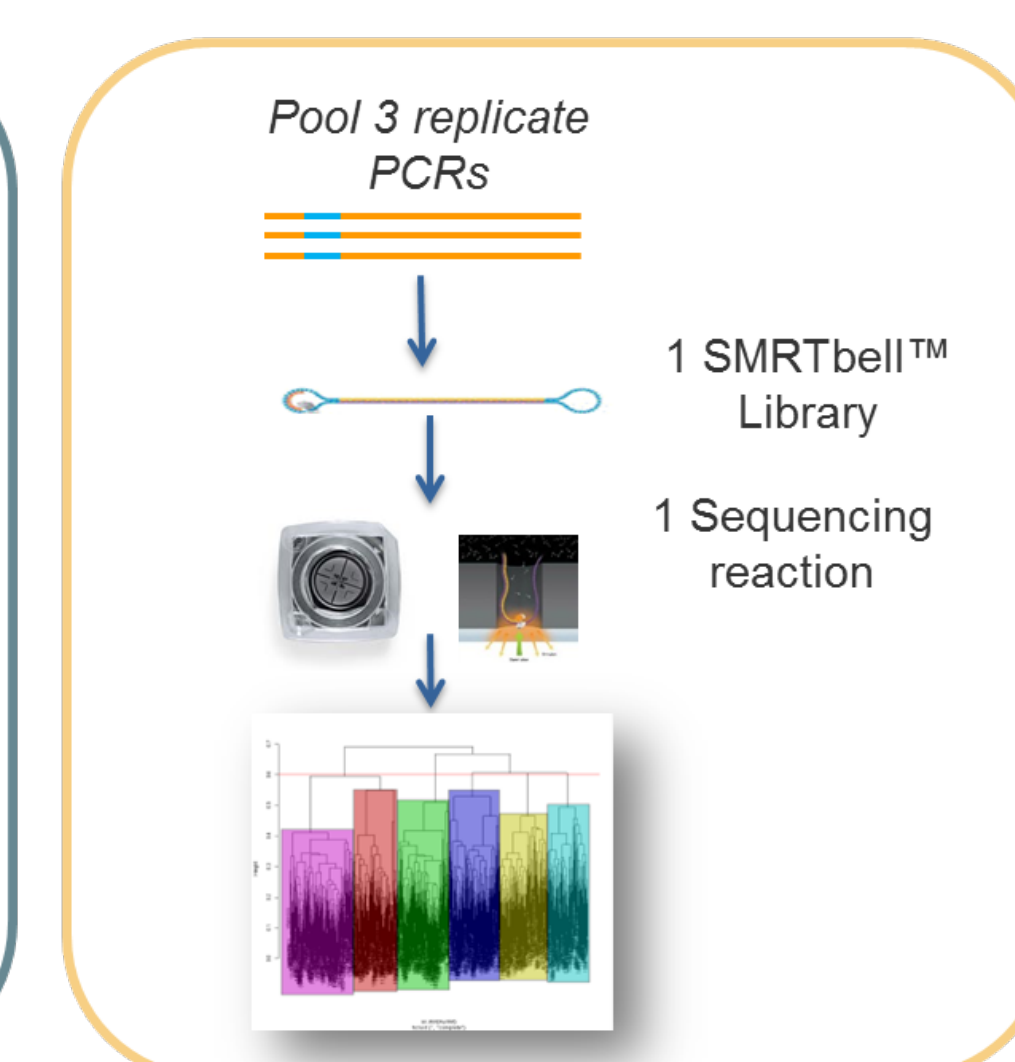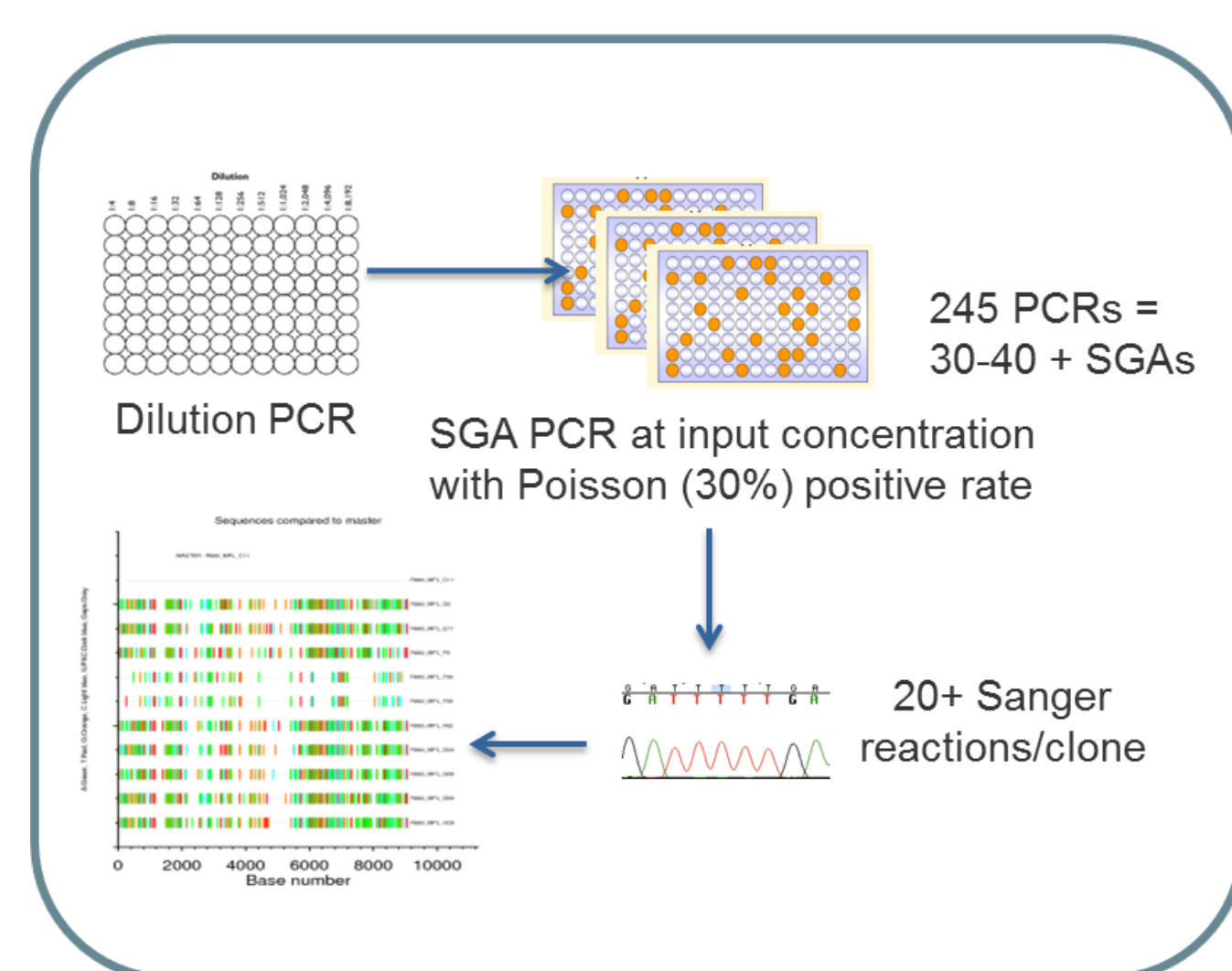## SMRT® Sequencing of Intact HIV-1 Genomes



**Full Genome – 9,084 bp**

## Two Methods to Measure the Genetic Diversity of Viral Infection



**Single Genome Amplification**
**30 viral genomes**
**8 – 10 days**

**PacBio® Sequencing**
**~1000 – 10,000 viral genomes**
**< 1 day**

## Mixtures of HIV-1: One or Two Base Differences

- Sequenced a sample containing a uniform mixture of HIV-1 clones differing by only one or two bases
- Analyzed using CluCon software
- **RESULT**: Five genomes estimated in the sample with 100% consensus accuracies across the entire ~9 kb genomes

```
Estimated Genome Haplotypes:

0 'AAAT+AAAAAAG+GT+GCA+ATTTACC+ACCC+' 25.2%

1 'AAT +AAAAAAG+AT+GGA+ATTTACC+ACCC+' 21.6%

2 'AAAT+AAAAAAG+AT+GGA+ATTTACC+ACCC+' 21.0%

3 'AAAT+AG AAAG+AT+GGA+ATTTACC+ACCC+' 19.6%

4 'AAAT+AAAAAAG+AT+GGA+ATTTTAC A CCC+' 12.6%
```

## Near-Full-Length Genomes from Mixture



- 100% correct near-full-length HIV-1 genomes with differences of five strains in the mixture highlighted by different colors
  - Two genomes differ by one base
  - Two genomes differ by two bases
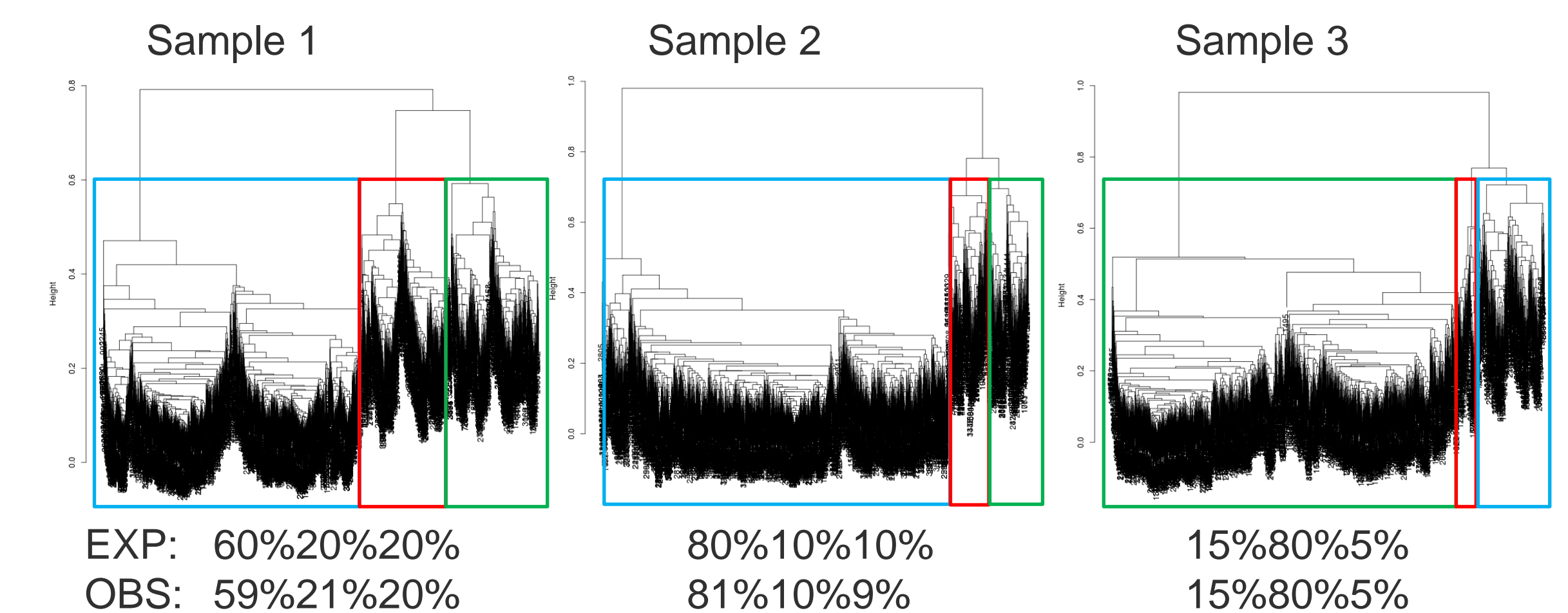  - Single genome had no differences

## CluCon Fine Clustering Technique

- Consensus Linear System Basis Function Deconvolution
- Identical to CluCon Strong Clustering but called when the number of minor variant column regions is too small for binomial bounds to have power
- Tally the read identities at the variant sites, the putative haplotypes counts
- Estimate true haplotypes from the observed reads by discounting noisy haplotypes using basis functions representing sequencing noise ( point spread functions in signal processing)

| Observed Frequency | Haplotypes | Basis Spread | Estimated Frequency |
|---|---|---|---|
| | H1: C+T+AAAA | Noise | |
| | H2: C+T+AAAAA | Distinct | |
| | H3: C+G+AAAA | Noise | |
| | H4: C+G+AAAAA | | |
| … | … | | … |

## Mixture of HIV-1 Genomes

- Sequenced three samples containing synthetic mixtures of HIV-1 clones at different abundances
- **RESULT**: Three genomes estimated in each sample with 100% consensus accuracies across the entire ~9 kb genomes and within 1% of expected abundances



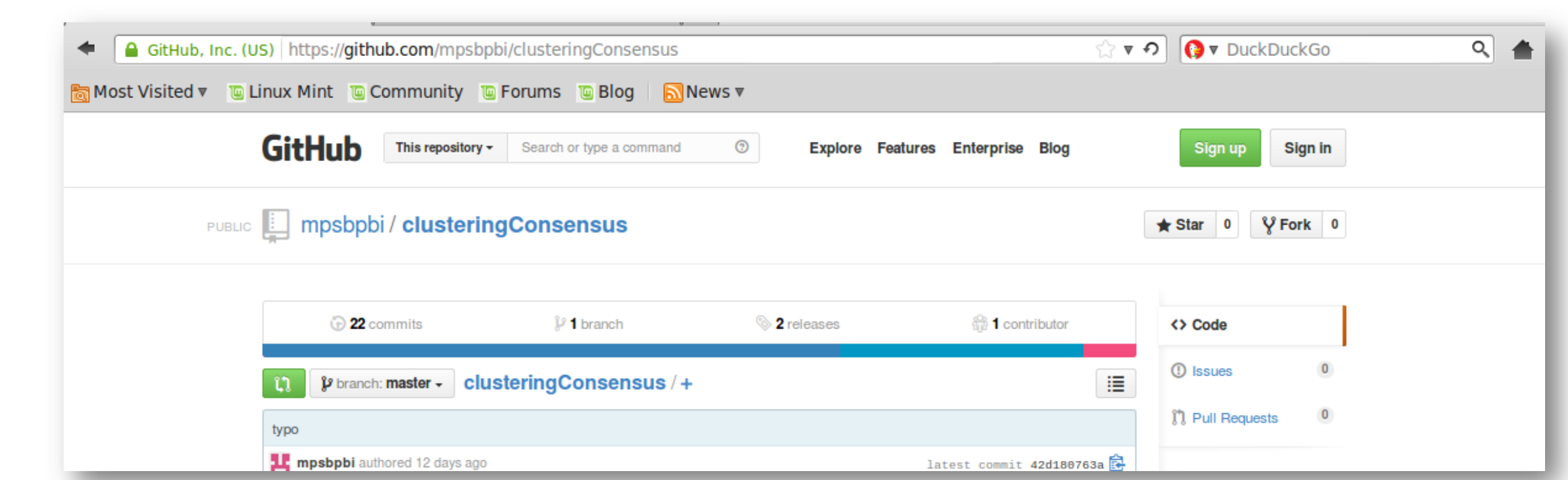| | Sample 1 | Sample 2 | Sample 3 |
|---|---|---|---|
| EXP: | 60%20%20% | 80%10%10% | 15%80%5% |
| OBS: | 59%21%20% | 81%10%9% | 15%80%5% |

## CluCon Strong Clustering Technique

- Consensus Sequence Clustering Under Binomial Bounds:
  - Generate alignment to run-specific consensus
  - Identify minor variant column regions
  - If no minor variants, return consensus
  - If divergence between read pairs is larger than expected by noise, then separate out as subcluster.
  - Recursively work on subclusters
- Clones in mixtures are at ~5% divergence (or 500 positions)
- Mathematically almost impossible to confuse one clone for another under sequencing with these divergences

## CluCon Software

- CluCon software and data available:

  https://github.com/mpsbpbi/clusteringConsensus



## Conclusions

- Complete characterization of HIV-1 genomes from single molecules
  - Sanger-quality, fully phased across entire genome
  - One SMRT Cell of sequencing
  - From samples with possibly complex mixtures of genomes even differing by single bases
- Sequencing and methods useful in other applications
  - Viral
  - Cancer
  - Metagenomics